

Quantifying Information Leakage in Document Redaction

Daniel Lopresti

Computer Science & Engineering
Lehigh University
Bethlehem, PA 18015, USA
lopresti@cse.lehigh.edu

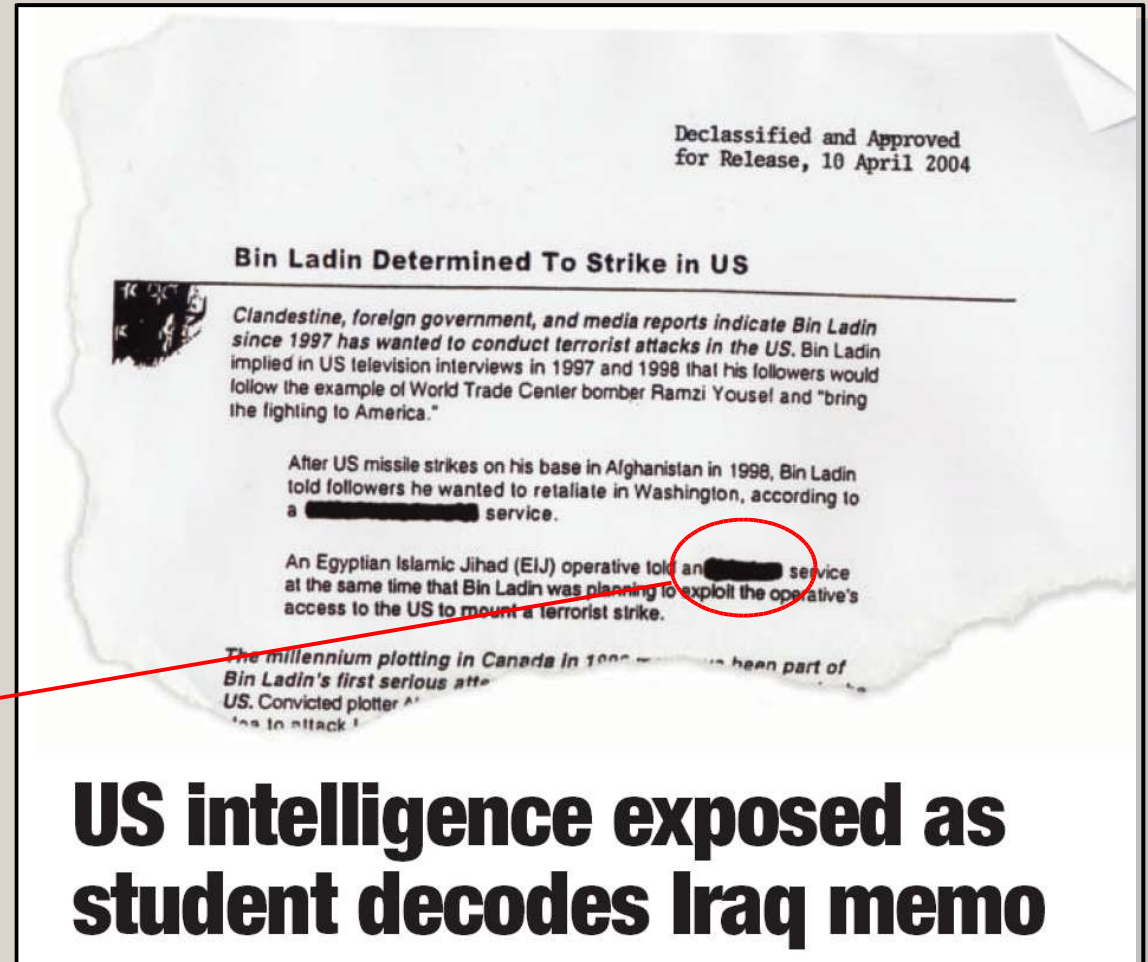
A. Lawrence Spitz

DocRec Ltd
34 Strathaven Place
Atawhai, Nelson, New Zealand
spitz@docrec.com

Motivation

Work by computer security researchers David Naccache and Claire Whelan as reported in *Nature*, May 2004.

“Egyptian” =



D. Butler, "US intelligence exposed as student decodes Iraq memo," *Nature*, 429:116, May 2004.

Is this a topic for research?

Brings together known techniques from document analysis and natural language processing in novel, perhaps interesting ways.

Some issues:

- No access to real pre-redacted data (of course) – it's confidential. Instead, must make assumptions and hope they're reasonable.
- E.g., leaks are unintended, not strategic (but that's also interesting).
- Optimization problem – minimal redaction needed to declassify.
- Attack need not be fully automated – semi-automated is sufficient.

Ultimate goals:

- Develop understanding of how (and how much) information leaks.
- Design PASS / FAIL test for deciding if OK to release document.

How might information leak?

- Text not completely obliterated. E.g., reflective qualities of “black” may differ for laserprinter toner and marker pen.
- While obscured, certain features still deducible. E.g., numbers and locations of ascender and descender characters.
- Exploiting string set-width in monospaced fonts (e.g., Courier). Combined with language modeling techniques, this can reveal missing text or at least limit possibilities.
- Exploiting string set-width in proportionally-spaced fonts (e.g., Times). Surprisingly, this reveals even more information ...

Tools for mounting attacks

- Image processing. Apply same sorts of techniques we already use in document analysis. E.g., histograms and adaptive thresholding.
- Font metrics. Many documents prepared using one of a few, well-known fonts. Font metric data is easily available (e.g., Adobe Font Metrics files). Naccache and Whelan first did font ID via simple image processing, then applied language modeling.
- Artifacts. Knowledge of ascenders, descenders, and i-dots may be sufficient to apply Character Shape Coding (Spitz, et al.).
- Natural Language Processing. Applicable in all of above cases. Public domain tools exist for text processing (e.g., tokenization, part-of-speech tagging). Internet makes building lexicons easy.

A. L. Spitz, "Progress in document reconstruction," *International Conference on Pattern Recognition*, pp. 464-467, 2002.



Simple image processing attack

Exploit differences in reflective qualities of “black.”

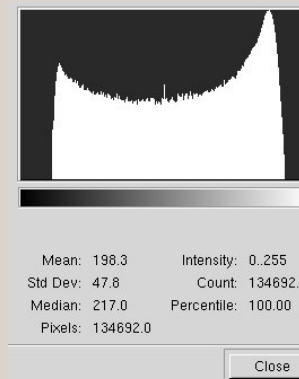
Obscure with black marker pen, then photocopy ...

he individual investor than Charles Schwab. We're consta
daries to better serve your needs. And today, I'm proud t
ver with the introduction of [REDACTED] Personal Choice™—
control over the way you invest in your hands, not your

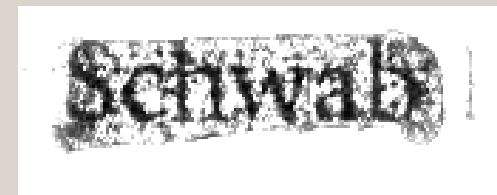
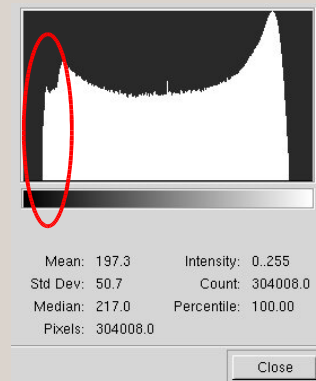
Threshold for redaction ...

he individual investor than Charles Schwab. We're consta
daries to better serve your needs. And today, I'm proud t
ver with the introduction of Schwab Personal Choice™—
control over the way you invest in your hands, not your

*Histogram
w/o redaction*



*Histogram
w/ redaction*



What does string set-width reveal?

*= width of two spaces
+ Senator's name*

Senator [REDACTED] of [REDACTED] introduced
bill under consideration on the floor of the United States

*Likely name
of US Senator*

*Probably
his/her state*

Preliminary experiments:

- Collect sample lexicons and font metric data.
- Study range of possible string set-widths and what it tells us.

Sample lexicons (all public domain)

YAWL

- “Yet another word list” is list of over 264,000 English words.

COUNTRIES

- 416 country names from around world (official and informal).

CONGRESS

- Names of 101 Senators (including VP) and 439 Representatives currently serving in U.S. Congress.

NAMES

- Cross-product of two lists from U.S. Census Bureau. First is list of male (1,219) and female (4,275) first names, while second is list of last names (88,799). Total of 487,861,706 names are generated.

Preliminary evaluation 1

Take Adobe Font Metrics files for Times, Helvetica, and Courier and count average number of strings of given set-width:

<i>Lexicon (Size)</i>	<i>Font</i>		
	<i>Times</i>	<i>Helvetica</i>	<i>Courier</i>
YAWL (264,057)	290	261	548
COUNTRIES (416)	1.33	1.28	1.97
CONGRESS (540)	1.66	1.60	2.90
NAMES (487,861,706)	481,125	427,573	969,903

Conclusions:

- With reliable small lexicon, attacker nearly always succeeds.
- Courier is “safer” font than Times or Helvetica (counter-intuitive).

Preliminary evaluation 2

Instead of average case, now look at worst-case analysis – number of strings which share same width with at most two other text strings:

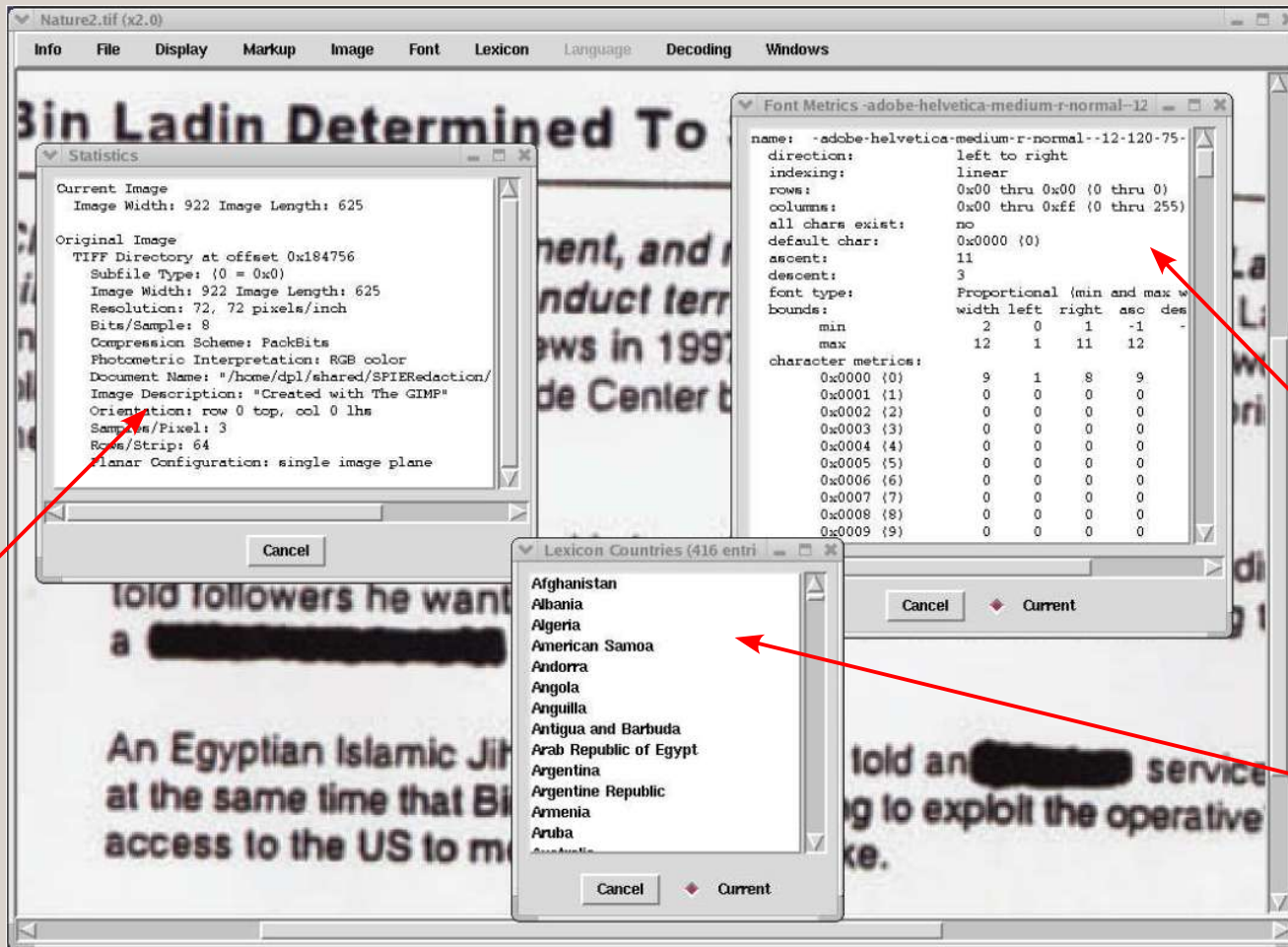
<i>Lexicon (Size)</i>	<i>Font</i>		
	<i>Times</i>	<i>Helvetica</i>	<i>Courier</i>
YAWL (264,057)	213	234	112
COUNTRIES (416)	394	393	260
CONGRESS (540)	472	514	233
NAMES (487,861,706)	49	78	20

Conclusion: even for very large lexicons, some strings are easily exposed using such techniques.

Plumber: a tool for finding leaks

- As previously noted, we strongly believe that semi-automated approaches lead to the most effective attacks on redacted text.
- We have implemented a prototype system to test some of these ideas. *Plumber* is written in Tcl/Tk, a popular scripting language for building applications with rich graphical user interfaces.
- *Plumber* is based around an image browser. User interacts with page in question, marking it up to delineate regions of interest, designate suspected font, and choose lexicon resources.
- *Plumber* then proposes possible interpretations for redacted region. Can render candidate strings in indicated font to overlay of page image to confirm guesses.
- Also implements wild-card search using character shape codes.

Plumber screen snapshot



Page image statistics

Adobe Font Metrics

Current lexicon

Exploiting string set-width

ill provide developments on the code-a-phone and w

Estimates for space and redaction widths

Candidate text overlay

(14,43) (301,57)

on Senator **Arlen Specter** o

ion bill under consideration

fined benefit pension plan a

ed in number and content ac

Estimated String Widths #0

- +10: Mike Dewine (14 + 283 + 14)
- +10: John F. Kerry (14 + 283 + 14)
- +10: Jeff Sessions (14 + 283 + 14)
- +14: Conrad Burns (14 + 287 + 14)
- +14: Gordon Smith (14 + 287 + 14)
- +18: Wayne Allard (14 + 291 + 14)
- +18: Thad Cochran (14 + 291 + 14)
- +18: Arlen Specter (14 + 291 + 14)
- +18: Craig Thomas (14 + 291 + 14)
- +22: Larry E. Craig (14 + 295 + 14)
- +22: John Edwards (14 + 295 + 14)

U.S. Senators (101 entries)

Image: scan1red.tif

Stated Image Resolution: 300 x 300 pixels/inch

Font: -adobe-times-medium-r-normal--12-120-75-75-p-64-iso8859-1

Image Resolution: 300

Space Width: 14

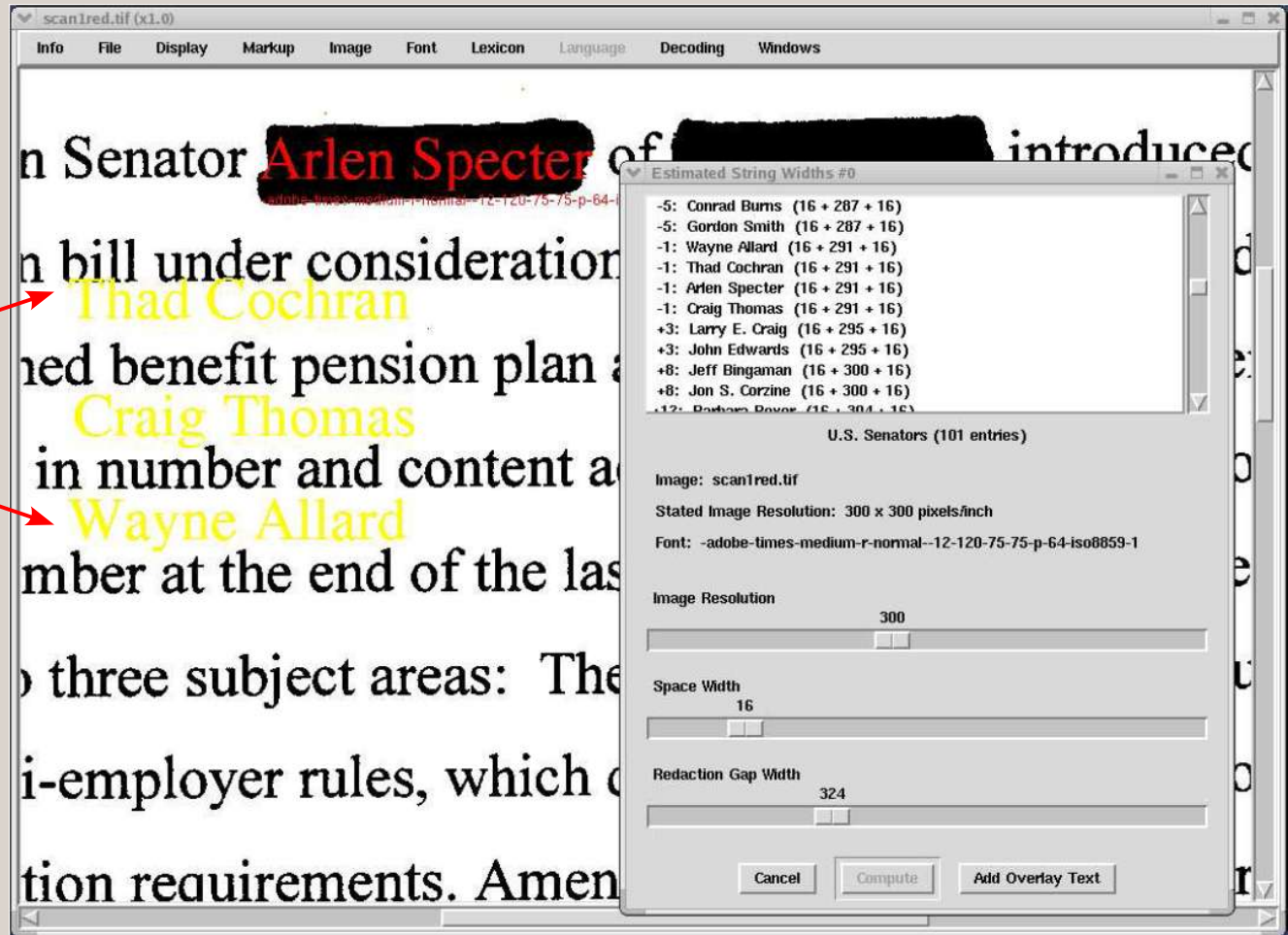
Redaction Gap Width: 301

Cancel Compute Add Overlay Text

Ranked list of candidate strings

Confirming candidate strings

Renderings
of alternate
candidates

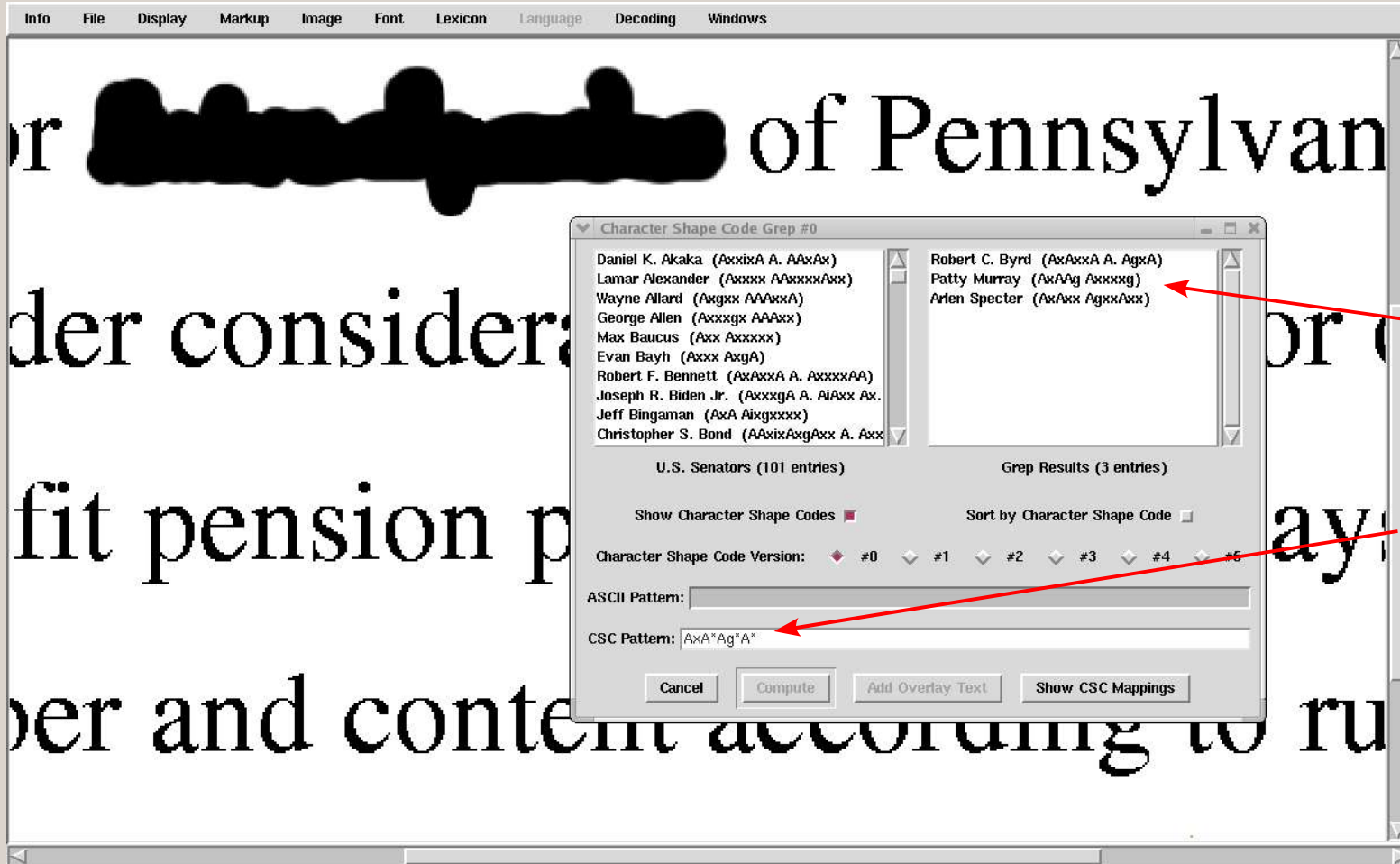


Character shape code mappings

Characters	V ₀	V ₁	V ₂	V ₃	V ₄	V ₅
amorsuvxwz	x	x	x	x	x	x
n			n			
c		e			c	
e					e	
ACGIOQSTUVWXYZflt		A		A	A	A
HMN				N		N
bhkL				b		b
BDEKR				E		E
PF				P		P
dJ				d		d
gpqy				g		
i				i		
j				j		

A. L. Spitz, "Using character shape codes for word spotting in document images," *Shape, Structure and Pattern Recognition*, pp. 382-389, 1995.

Plumber CSC wild-card search



Matching candidates

CSC pattern = Ax*A*Ag*A*

Plumber CSC wild-card search

A Senate USOC Reform Bill (No. S.1404 in the
[redacted] of [redacted] was on the verge of being
before a vote on the

Character Shape Code Mappings

Characters	V ₀	V ₁	V ₂	V ₃	V ₄	V ₅
amorsuvvwz	x	x	x	x	x	x
n			n			
c		e			c	
e						e
ACGIOQSTUVWXYZflt		A		A	A	A
HMN			N	N	N	
bhKL			b	b	b	
BDEKR			E	E	E	
PF			P	P	P	
dJ			d	d	d	
gpqy				g		
i				i		
j					j	

Character Shape Code Grep #1

Alabama (AAxAxxx)
Alaska (AAxxAx)
American Samoa (Axxxxxxx Axxxx)
Arizona (Axixxxx)
Arkansas (AxAxxxxx)
California (AxAlAxxxxix)
Colorado (AxAxxxAx)
Connecticut (AxxxxxAxxxA)
Delaware (AxAvxxxx)
District of Columbia (AxAxixA xA AxAvxxA)

Arizona (Axixxxx)
Maine (Axixx)

U.S. States (59 entries) Grep Results (2 entries)

Show Character Shape Codes Sort by Character Shape Code

Character Shape Code Version: #0 #1 #2 #3 #4 #5

ASCII Pattern:

CSC Pattern: Axix*

Cancel Compute Add Overlay Text Show CSC Mappings

Only 2 states
match: Arizona
and Maine

CSC pattern
= Axix*

Related (but even harder) problem

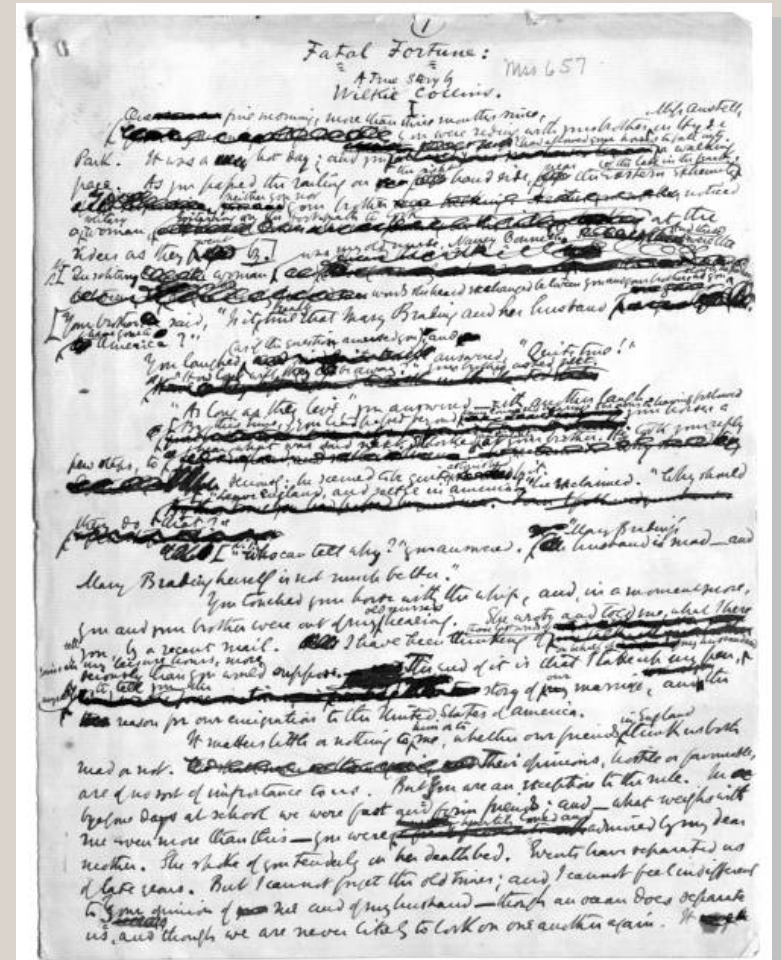
Interestingly, this problem does not appear to be unique to redacted text.

Title page from “Fatal Fortune: A True Story” by Wilkie Collins.

“Collins' manuscript is a short story told from the perspective of a woman who falls in love with a man suspected of being 'mad.' He has been subsequently disenfranchised of his fortune by scheming executors.”

From “I remain: A Digital Archive of Letters, Manuscripts, and Ephemera” at Lehigh University.

<http://digital.lib.lehigh.edu/remain/>



Summary

Whether the problem of information leakage through redaction is of practical importance is unclear to us (and might forever remain so).

Still, there seem to be many interesting technical issues here:

- How can we quantify the amount of information leakage?
- When is a document safe for release?
- Can the process of redaction (or confirming that a document has been safely redacted) be automated?
- Is it possible to design effective counter-measures to make documents safer? E.g., special fonts or typesetting conventions that defeat set-width attacks?
- Does what we learn from studying this problem apply elsewhere?